



北京理工大学
BEIJING INSTITUTE OF TECHNOLOGY

大数据泛构实验报告 3

——树状度量空间索引实现与性能对比

学 院:	珠海校区
专 业:	计算机科学与技术
学 号:	3120256739
姓 名:	丁纪翔
任课教师:	毛睿

2025 年 12 月 10 日

目录

1. 引言	5
1.1 研究背景与意义	5
1.2 任务回顾与目标	5
1.2.1 Assignment 1 & 2 工作简述	5
1.2.2 Assignment 3 核心目标	5
1.3 实验环境	6
1.4 报告结构	6
2. 树状度量空间索引理论基础	7
2.1 树状索引概述	7
2.1.1 基于 Pivot Table 的索引回顾	7
2.1.2 树状索引的基本思想	7
2.1.3 数据划分方式分类	7
2.2 GH 树 (Generalized Hyperplane Tree)	8
2.2.1 基本思想：超平面划分	8
2.2.2 数据结构设计	8
2.2.3 批建算法原理	9
2.2.4 范围查询算法原理	9
2.3 VP 树 (Vantage Point Tree)	10
2.3.1 基本思想：球形划分	10
2.3.2 数据结构设计	10
2.3.3 批建算法原理	11
2.3.4 范围查询算法原理	11
2.4 GHT 与 VPT 的理论对比	12
2.4.1 数据划分方式对比	12
2.4.2 支撑点使用对比	12
2.4.3 空间划分特性对比	12
3. GHT 和 VPT 实现	13
3.1 系统架构扩展	13
3.1.1 在原有架构中集成树状索引模块	13
3.1.2 更新后的系统模块划分	13

3.2	GHT 实现	13
3.2.1	GHT 节点数据结构实现	13
3.2.2	GHT 批建算法实现	14
3.2.3	GHT 范围查询实现	16
3.2.4	GHT kNN 查询实现	17
3.3	VPT 实现	17
3.3.1	VPT 节点数据结构实现	17
3.3.2	VPT 批建算法实现	18
3.3.3	VPT 范围查询实现	20
3.3.4	VPT kNN 查询实现	21
3.4	Pivot 选择与树高控制	21
3.4.1	统一的 Pivot 选择策略	21
3.4.2	树高控制方法	22
4.	功能正确性验证	23
4.1	测试环境与小规模测试数据集	23
4.2	GHT 正确性验证	23
4.2.1	构建过程展示与分析	23
4.2.2	查询过程展示与手工验证	25
4.2.3	与线性扫描结果对比	25
4.3	VPT 正确性验证	26
4.3.1	构建过程展示与分析	26
4.3.2	查询过程展示与手工验证	27
4.3.3	与线性扫描结果对比	28
4.4	GHT 与 VPT 结果一致性验证	28
4.4.1	范围查询一致性测试	28
4.4.2	kNN 查询一致性测试	28
4.4.3	不同 Pivot 策略下的一致性	28
4.5	单元测试结果汇总	29
5.	性能对比实验设计与实施	29
5.1	实验方案设计	29
5.1.1	性能评价指标设计及理由	29

5.1.2	数据集选择及说明	30
5.1.3	查询类型、查询对象与参数设置	30
5.1.4	树高控制方法说明	31
5.1.5	GHT 和 VPT 的统一参数设置	31
5.2	蛋白质序列数据集实验	31
5.2.1	数据集特征分析	31
5.2.2	索引构建性能对比	32
5.2.3	范围查询性能对比	32
5.2.4	kNN 查询性能对比	32
5.3	低维向量数据集实验	32
5.3.1	数据集特征分析	32
5.3.2	索引构建性能对比	33
5.3.3	范围查询性能对比	33
5.3.4	kNN 查询性能对比	33
5.4	高维向量数据集实验	33
5.4.1	数据集特征分析	33
5.4.2	索引构建性能对比	34
5.4.3	范围查询性能对比	34
5.4.4	kNN 查询性能对比	34
5.5	参数影响分析	34
5.5.1	最大叶子节点大小的影响	34
5.5.2	Pivot 选择策略的影响	35
6.	性能对比分析	35
6.1	实验结果汇总	35
6.2	不同数据集上的性能差异分析	36
6.2.1	蛋白质序列数据分析	36
6.2.2	低维向量数据分析	36
6.2.3	高维向量数据分析	37
6.3	性能差异的原因分析	37
6.3.1	数据划分方式的影响	37
6.3.2	支撑点使用效率的影响	38
6.3.3	数据分布特征的影响	38

6.3.4	维度灾难的影响	38
6.4	GHT 与 VPT 的优缺点总结	39
6.4.1	GHT 的优势与局限	39
6.4.2	VPT 的优势与局限	39
6.4.3	适用场景分析	40
6.5	改进方向讨论	40
6.5.1	Pivot 选择策略优化	40
6.5.2	树结构平衡性优化	40
6.5.3	多路划分扩展	41
6.5.4	动态索引维护	41
6.5.5	高维数据优化	41
7.	总结与展望	41
7.1	工作总结	41
7.1.1	代码实现	41
7.1.2	正确性验证	42
7.1.3	性能实验	42
7.1.4	分析与讨论	42
7.2	主要结论	43
7.3	未来展望	43
7.3.1	算法优化	43
7.3.2	扩展应用	43
7.3.3	性能提升	44
7.3.4	理论分析	44
7.4	结语	44

1. 引言

1.1 研究背景与意义

度量空间 (Metric Space) 是一种通用的数据抽象方式, 可以涵盖向量、字符串、图像、生物序列等多种数据类型。在度量空间中, 数据对象之间的相似性通过满足正定性、对称性和三角不等式的距离函数来度量。相似性查询是度量空间数据管理中最基本的操作之一, 包括范围查询 (Range Query) 和 k 近邻查询 (kNN Query)。

随着大数据时代的到来, 如何高效地处理海量度量空间数据成为重要的研究课题。线性扫描方法虽然简单, 但时间复杂度为 $O(n)$, 无法满足大规模数据集的查询需求。因此, 设计高效的索引结构来加速相似性查询具有重要的理论意义和实际价值。

树状索引是度量空间索引的重要类别, 通过层次化的数据划分实现高效的剪枝, 可以显著减少查询时的距离计算次数。GH 树 (Generalized Hyperplane Tree) 和 VP 树 (Vantage Point Tree) 是两种经典的树状索引结构, 分别采用超平面划分和球形划分策略。

1.2 任务回顾与目标

1.2.1 Assignment 1 & 2 工作简述

在 Assignment 1 中, 我们建立了度量空间数据管理的基础设施, 包括:

- 核心抽象类: `MetricSpaceData` 和 `MetricFunction`
- 具体数据类型: 向量数据 (`VectorData`) 和蛋白质序列 (`ProteinData`)
- 距离函数实现: 闵可夫斯基距离和 `Alignment` 距离
- 数据读取模块: 支持从文件读取向量和蛋白质序列数据

在 Assignment 2 中, 我们实现了基础的查询算法和 Pivot Table 索引:

- 线性扫描查询: 范围查询、kNN 查询、多样化 kNN 查询
- Pivot Table 索引: 利用支撑点预计算距离实现查询剪枝
- Pivot 选择策略: 随机选择、FFT、增量选择等

1.2.2 Assignment 3 核心目标

本次 Assignment 3 的核心目标是实现两种树状度量空间索引——GH 树和 VP 树, 并进行全面的性能对比分析。具体任务包括:

1. 代码实现：实现 GH 树和 VP 树的数据结构、批建算法、范围查询和 kNN 查询
2. 实验设计：设计科学的性能对比实验方案，包括评价指标、数据集选择、参数设置等
3. 正确性验证：通过与线性扫描对比、手工验证等方式证明实现的正确性
4. 性能对比：在多种数据集上对比 GH 树和 VP 树的性能
5. 分析讨论：分析性能差异的原因，总结两种索引的优缺点，讨论改进方向

1.3 实验环境

本实验的软硬件环境如表1-1所示。

表 1-1 实验环境配置

项目	配置
操作系统	Windows 11
CPU	Intel Core
内存	16GB
编程语言	Java 11
构建工具	Maven 3.8+
IDE	VS Code

1.4 报告结构

本报告的组织结构如下：

- 第 2 章树状度量空间索引理论基础：介绍 GH 树和 VP 树的基本概念、数据结构和算法原理
- 第 3 章 GHT 和 VPT 实现：详细描述系统架构扩展和核心算法的实现
- 第 4 章功能正确性验证：通过小规模数据演示和与线性扫描对比验证实现的正确性
- 第 5 章性能对比实验设计与实施：设计实验方案并在多种数据集上进行实验
- 第 6 章性能对比分析：分析实验结果，讨论性能差异的原因和改进方向
- 第 7 章总结与展望：总结本次工作，展望未来研究方向

本项目的完整代码已上传至 GitHub 仓库：

https://github.com/sylvanding/BigDataGenhierarchy_Jixiang_20251116

2. 树状度量空间索引理论基础

2.1 树状索引概述

2.1.1 基于 Pivot Table 的索引回顾

在 Assignment 2 中，我们实现了 Pivot Table 索引。Pivot Table 通过预计算数据到支撑点（Pivot）的距离，利用三角不等式进行剪枝。其基本原理如下：

对于查询对象 q 、数据对象 x 和支撑点 p ，由三角不等式可得：

$$|d(q, p) - d(x, p)| \leq d(q, x) \leq d(q, p) + d(x, p) \quad (1)$$

因此，如果 $|d(q, p) - d(x, p)| > r$ ，则 $d(q, x) > r$ ，可以排除 x 。

Pivot Table 的局限性包括：

- 每次查询需要计算到所有 pivot 的距离
- 只能利用全局的 pivot 信息，不能根据数据分布进行层次化剪枝
- 空间开销大，需要存储所有数据到所有 pivot 的距离

2.1.2 树状索引的基本思想

树状索引通过递归划分数据集，构建层次化的索引结构。其优势包括：

- 层次化剪枝：通过树的层次结构，可以提前排除整个子树
- 局部优化：每个节点可以根据子树的数据分布选择 pivot
- 距离计算少：只计算到访问路径上节点的 pivot 的距离
- 可扩展性好：适合大规模数据集

2.1.3 数据划分方式分类

根据划分方式，树状索引主要分为两类：

1. 超平面划分（**Hyperplane Partitioning**）：代表为 GH 树，使用两个支撑点定义超平面
2. 球形划分（**Ball Partitioning**）：代表为 VP 树，使用一个支撑点定义球形边界

图2-1展示了两种划分方式的几何直观。

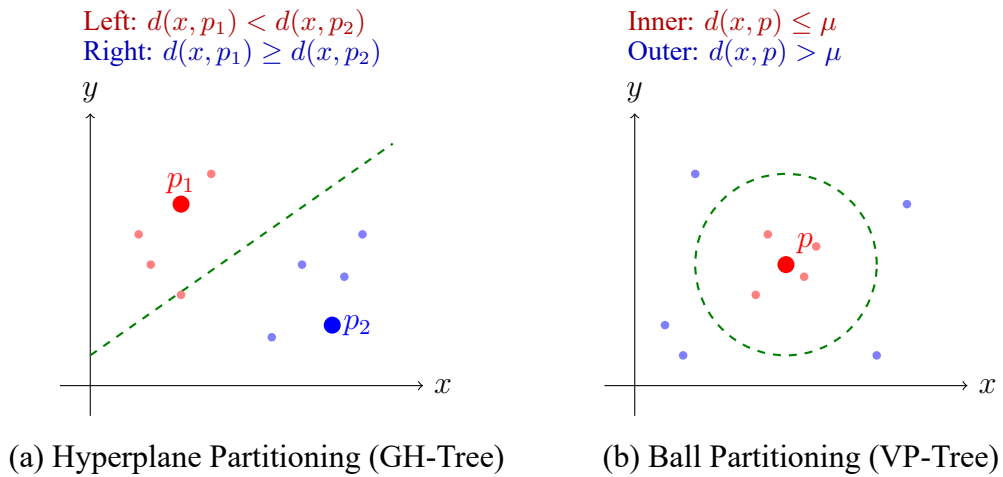


图 2-1 两种数据划分方式的几何直观

2.2 GH 树 (Generalized Hyperplane Tree)

2.2.1 基本思想：超平面划分

GH 树使用广义超平面来划分度量空间。给定两个支撑点 p_1 和 p_2 ，定义广义超平面为：

$$H(p_1, p_2) = \{x \in M \mid d(x, p_1) = d(x, p_2)\} \quad (2)$$

这个超平面将空间划分为两个区域：

- 左区域: $L = \{x \in M \mid d(x, p_1) < d(x, p_2)\}$
- 右区域: $R = \{x \in M \mid d(x, p_1) \geq d(x, p_2)\}$

在欧几里得空间中，满足 $d(x, p_1) = d(x, p_2)$ 的点集形成垂直平分线/面。在一般度量空间中，这个集合可能不是平面，故称“广义”超平面。

2.2.2 数据结构设计

GH 树的节点结构如下：

内部节点包含：

- 两个支撑点 p_1 和 p_2
- 左子树指针（离 p_1 更近的数据）
- 右子树指针（离 p_2 更近的数据）
- 节点深度

叶子节点包含：

- 数据对象列表
- 节点深度

2.2.3 批建算法原理

GH 树的批建算法采用递归方式，如算法1所示。

Algorithm 1 GH 树批建算法

Input: 数据集 S ，当前深度 $depth$

Output: GH 树节点

```

1: if  $|S| \leq \text{maxLeafSize}$  and  $depth \geq \text{minTreeHeight}$  then
2:   return  $\text{LeafNode}(S, depth)$ 
3: end if
4:  $(p_1, p_2) \leftarrow \text{SelectTwoPivots}(S)$ 
5:  $S_{\text{left}} \leftarrow \{x \in S \mid d(x, p_1) < d(x, p_2)\}$ 
6:  $S_{\text{right}} \leftarrow \{x \in S \mid d(x, p_1) \geq d(x, p_2)\}$ 
7:  $\text{leftChild} \leftarrow \text{BuildGHTree}(S_{\text{left}}, depth + 1)$ 
8:  $\text{rightChild} \leftarrow \text{BuildGHTree}(S_{\text{right}}, depth + 1)$ 
9: return  $\text{GHInternalNode}(p_1, p_2, \text{leftChild}, \text{rightChild}, depth)$ 

```

2.2.4 范围查询算法原理

GH 树的剪枝规则是其查询效率的关键。设查询对象为 q ，查询半径为 r ， $d_1 = d(q, p_1)$ ， $d_2 = d(q, p_2)$ 。

剪枝规则 1（排除左子树）：如果 $d_1 - d_2 > 2r$ ，则可以排除左子树。

证明：左子树中任意数据点 x 满足 $d(x, p_1) < d(x, p_2)$ 。由三角不等式：

$$d(q, x) \geq d(q, p_1) - d(x, p_1) = d_1 - d(x, p_1) \quad (3)$$

$$d(q, x) \geq d(x, p_2) - d(q, p_2) = d(x, p_2) - d_2 \quad (4)$$

两式相加： $2d(q, x) \geq d_1 - d(x, p_1) + d(x, p_2) - d_2 > d_1 - d_2$

若 $d_1 - d_2 > 2r$ ，则 $d(q, x) > r$ ，可排除。

剪枝规则 2（排除右子树）：如果 $d_2 - d_1 > 2r$ ，则可以排除右子树。

范围查询算法如算法2所示。

Algorithm 2 GH 树范围查询算法

Input: 查询对象 q , 查询半径 r , 当前节点 $node$

Output: 满足条件的数据对象列表

```

1: if  $node$  is LeafNode then
2:    $result \leftarrow \emptyset$ 
3:   for each  $x$  in  $node.data$  do
4:     if  $d(q, x) \leq r$  then
5:       add  $x$  to  $result$ 
6:     end if
7:   end for
8:   return  $result$ 
9: end if
10:  $result \leftarrow \emptyset$ 
11:  $d_1 \leftarrow d(q, node.p_1)$ 
12:  $d_2 \leftarrow d(q, node.p_2)$ 
13: if NOT  $(d_1 - d_2 > 2r)$  then
14:    $result.addAll(GHRangeQuery(q, r, node.leftChild))$ 
15: end if
16: if NOT  $(d_2 - d_1 > 2r)$  then
17:    $result.addAll(GHRangeQuery(q, r, node.rightChild))$ 
18: end if
19: return  $result$ 

```

2.3 VP 树 (Vantage Point Tree)

2.3.1 基本思想：球形划分

VP 树使用球形区域来划分度量空间。给定一个支撑点 (Vantage Point) p 和划分半径 μ , 空间被划分为:

- 内球: $B_{inner} = \{x \in M \mid d(x, p) \leq \mu\}$
- 外球: $B_{outer} = \{x \in M \mid d(x, p) > \mu\}$

通常 μ 取所有数据到 p 距离的中位数, 以保证划分平衡。

2.3.2 数据结构设计

VP 树的节点结构与 GH 树类似, 但有重要区别:

内部节点包含:

- 一个支撑点 p
- 划分中位数 μ
- 内球子树 (距离范围 $[L_{inner}, U_{inner}]$)

- 外球子树（距离范围 $[L_{outer}, U_{outer}]$ ）
- 节点深度

关键信息是每个子树的距离范围 $[L, U]$ ，表示该子树中数据到 **pivot** 的距离的最小值和最大值。

2.3.3 批建算法原理

VP 树的批建算法如算法3所示。

Algorithm 3 VP 树批建算法

Input: 数据集 S ，当前深度 $depth$

Output: VP 树节点

```

1: if  $|S| \leq \text{maxLeafSize}$  and  $depth \geq \text{minTreeHeight}$  then
2:   return  $\text{LeafNode}(S, depth)$ 
3: end if
4:  $p \leftarrow \text{SelectVantagePoint}(S)$ 
5:  $\text{distances} \leftarrow \{(x, d(x, p)) \mid x \in S, x \neq p\}$ 
6: Sort  $\text{distances}$  by distance value
7:  $\text{mid} \leftarrow |\text{distances}|/2$ 
8:  $S_{\text{inner}} \leftarrow \text{distances}[0 : \text{mid}]$ 
9:  $S_{\text{outer}} \leftarrow \text{distances}[\text{mid} : \text{end}]$ 
10:  $\text{range}_{\text{inner}} \leftarrow [\min(S_{\text{inner}}.\text{distances}), \max(S_{\text{inner}}.\text{distances})]$ 
11:  $\text{range}_{\text{outer}} \leftarrow [\min(S_{\text{outer}}.\text{distances}), \max(S_{\text{outer}}.\text{distances})]$ 
12:  $\text{innerChild} \leftarrow \text{BuildVPtree}(S_{\text{inner}}, depth + 1)$ 
13:  $\text{outerChild} \leftarrow \text{BuildVPtree}(S_{\text{outer}}, depth + 1)$ 
14: return  $\text{VPInternalNode}(p, [\text{innerChild}, \text{outerChild}], [\text{range}_{\text{inner}}, \text{range}_{\text{outer}}], depth)$ 

```

2.3.4 范围查询算法原理

VP 树的剪枝基于距离范围信息。设查询对象为 q ，查询半径为 r ， $d_q = d(q, p)$ ，子树 i 的距离范围为 $[L_i, U_i]$ 。

剪枝规则 1（内侧剪枝）：如果 $d_q + r < L_i$ ，则可以排除子树 i 。

几何解释：查询球的最远点到 p 的距离为 $d_q + r$ ，如果小于子树中最近点的距离 L_i ，则不相交。

剪枝规则 2（外侧剪枝）：如果 $d_q - r > U_i$ ，则可以排除子树 i 。

几何解释：查询球的最近点到 p 的距离为 $d_q - r$ ，如果大于子树中最远点的距离 U_i ，则不相交。

图2-2展示了 VP 树的剪枝规则。

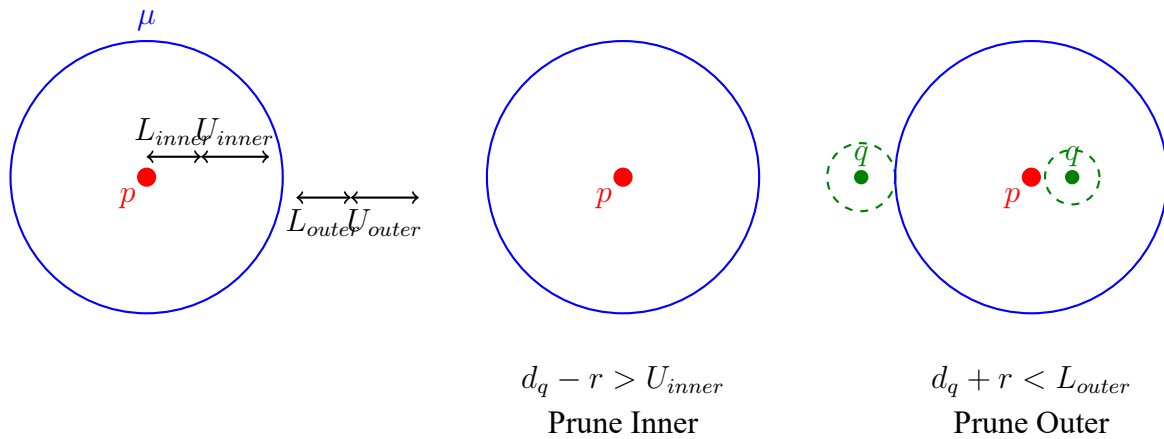


图 2-2 VP 树剪枝规则示意

2.4 GHT 与 VPT 的理论对比

2.4.1 数据划分方式对比

表2-2对比了 GH 树和 VP 树的划分方式。

表 2-2 GH 树与 VP 树数据划分方式对比

特性	GH 树	VP 树
划分方式	超平面划分	球形划分
几何形状	两点等距平面	同心球
划分平衡性	依赖数据分布	中位数保证平衡

2.4.2 支撑点使用对比

表2-3对比了两种树的支撑点使用情况。

表 2-3 支撑点使用对比

特性	GH 树	VP 树
每节点 pivot 数	2	1
pivot 作用	定义超平面	定义球心
距离计算（每节点）	2 次	1 次
剪枝信息	两距离差值	距离范围

2.4.3 空间划分特性对比

GH 树的特点：

- 优势：划分更对称，在低维空间效果好

- 劣势：需要 2 个 pivot，构建开销略大

VP 树的特点：

- 优势：只需 1 个 pivot，距离范围提供精确剪枝
- 劣势：球形划分可能导致空间浪费，高维情况受维度灾难影响

3. GHT 和 VPT 实现

3.1 系统架构扩展

3.1.1 在原有架构中集成树状索引模块

在 Assignment 1 和 2 的基础上，我们扩展系统架构以支持树状索引。图3-3展示了扩展后的系统模块关系。

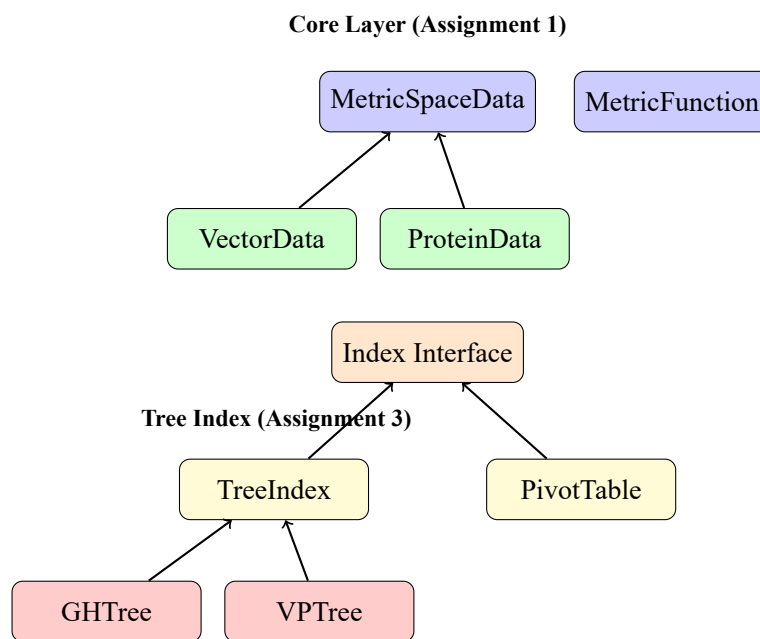


图 3-3 系统架构图

3.1.2 更新后的系统模块划分

表3-4展示了 Assignment 3 新增的主要模块。

3.2 GHT 实现

3.2.1 GHT 节点数据结构实现

GH 树内部节点的核心实现如下：

表 3-4 Assignment 3 新增模块

模块/类	功能说明
index.Index	所有索引结构的统一接口
index.tree.TreeIndex	树索引的抽象基类
index.tree.TreeNode	树节点接口
index.tree.InternalNode	内部节点抽象类
index.tree.LeafNode	叶子节点类
index.tree.common.TreeConfig	树配置类
index.tree.common.TreeHeightController	树高控制器
index.tree.ghTree.GHTree	GH 树实现
index.tree.ghTree.GHInternalNode	GH 树内部节点
index.tree.vptree.VPTree	VP 树实现
index.tree.vptree.VPInternalNode	VP 树内部节点

```

1 public class GHInternalNode extends InternalNode {
2     public GHInternalNode(MetricSpaceData pivot1,
3                           MetricSpaceData pivot2,
4                           TreeNode leftChild,
5                           TreeNode rightChild,
6                           int depth) {
7         this.pivots = List.of(pivot1, pivot2);
8         this.children = List.of(leftChild, rightChild);
9         this.depth = depth;
10    }
11
12    public MetricSpaceData getPivot1() { return pivots.get(0); }
13    public MetricSpaceData getPivot2() { return pivots.get(1); }
14    public TreeNode getLeftChild() { return children.get(0); }
15    public TreeNode getRightChild() { return children.get(1); }
16 }

```

Listing 1: GHInternalNode 核心代码

3.2.2 GHT 批建算法实现

GH 树批建的核心逻辑实现如下：

```

1 protected TreeNode buildTreeRecursive(
2     List<? extends MetricSpaceData> data, int depth) {

```

```
3      // Step 1: Check termination condition
4      if (heightController.canCreateLeaf(depth, data.size())) {
5          return new LeafNode(data, depth);
6      }
7
8      // Step 2: Select two pivots
9      MetricSpaceData[] pivots = selectTwoPivots(data);
10     MetricSpaceData pivot1 = pivots[0];
11     MetricSpaceData pivot2 = pivots[1];
12
13     // Step 3: Partition data
14     List<MetricSpaceData> leftData = new ArrayList<>();
15     List<MetricSpaceData> rightData = new ArrayList<>();
16
17     for (MetricSpaceData obj : data) {
18         double d1 = metric.getDistance(obj, pivot1);
19         double d2 = metric.getDistance(obj, pivot2);
20         buildDistanceComputations += 2;
21
22         if (d1 < d2) {
23             leftData.add(obj);
24         } else {
25             rightData.add(obj);
26         }
27     }
28
29     // Step 4: Recursive build
30     TreeNode leftChild = buildTreeRecursive(leftData, depth + 1);
31     TreeNode rightChild = buildTreeRecursive(rightData, depth + 1);
32
33     return new GHInternalNode(pivot1, pivot2,
34                               leftChild, rightChild, depth);
35 }
```

Listing 2: GH 树批建算法核心实现

3.2.3 GHT 范围查询实现

GH 树范围查询实现了超平面剪枝规则：

```

1 public List<MetricSpaceData> rangeQuery(
2     MetricSpaceData queryObject, double radius) {
3     List<MetricSpaceData> result = new ArrayList<>();
4     resetStatistics();
5     rangeQueryRecursive(root, queryObject, radius, result);
6     return result;
7 }
8
9 private void rangeQueryRecursive(TreeNode node,
10     MetricSpaceData q, double r, List<MetricSpaceData> result) {
11     nodeAccesses++;
12
13     if (node.isLeaf()) {
14         LeafNode leaf = (LeafNode) node;
15         for (MetricSpaceData obj : leaf.getData()) {
16             double dist = metric.getDistance(q, obj);
17             queryDistanceComputations++;
18             if (dist <= r) {
19                 result.add(obj);
20             }
21         }
22         return;
23     }
24
25     GHInternalNode internal = (GHInternalNode) node;
26     double d1 = metric.getDistance(q, internal.getPivot1());
27     double d2 = metric.getDistance(q, internal.getPivot2());
28     queryDistanceComputations += 2;
29
30     // Pruning rule 1: if  $d1 - d2 > 2r$ , prune left
31     if (!(d1 - d2 > 2 * r)) {
32         rangeQueryRecursive(internal.getLeftChild(), q, r, result);
33     }
34

```

```

35 // Pruning rule 2: if  $d_2 - d_1 > 2r$ , prune right
36 if (!( $d_2 - d_1 > 2 * r$ )) {
37     rangeQueryRecursive(internal.getRightChild(), q, r, result);
38 }
39 }

```

Listing 3: GH 树范围查询实现

3.2.4 GHT kNN 查询实现

kNN 查询使用优先队列维护当前 k 个最近邻：

```

1 public List<MetricSpaceData> knnQuery(
2     MetricSpaceData queryObject, int k) {
3     PriorityQueue<DataWithDistance> knnQueue =
4         new PriorityQueue<>((a, b) ->
5             Double.compare(b.distance, a.distance)); // max-heap
6
7     resetStatistics();
8     knnQueryRecursive(root, queryObject, k, knnQueue);
9
10    // Convert to result list
11    List<MetricSpaceData> result = new ArrayList<>();
12    while (!knnQueue.isEmpty()) {
13        result.add(0, knnQueue.poll().data);
14    }
15    return result;
16 }

```

Listing 4: GH 树 kNN 查询实现

3.3 VPT 实现

3.3.1 VPT 节点数据结构实现

VP 树内部节点存储距离范围信息：

```

1 public class VPInternalNode extends InternalNode {

```

```

2     private List<DistanceRange> distanceRanges;
3
4     public static class DistanceRange {
5         public double lower; // minimum distance
6         public double upper; // maximum distance
7
8         public DistanceRange(double lower, double upper) {
9             this.lower = lower;
10            this.upper = upper;
11        }
12    }
13
14    public VPInternalNode(MetricSpaceData pivot,
15                          List<TreeNode> children,
16                          List<DistanceRange> distanceRanges,
17                          int depth) {
18        this.pivots = Collections.singletonList(pivot);
19        this.children = new ArrayList<>(children);
20        this.distanceRanges = new ArrayList<>(distanceRanges);
21        this.depth = depth;
22    }
23
24    public MetricSpaceData getPivot() { return pivots.get(0); }
25    public DistanceRange getDistanceRange(int i) {
26        return distanceRanges.get(i);
27    }
28 }

```

Listing 5: VPInternalNode 核心代码

3.3.2 VPT 批建算法实现

VP 树批建基于中位数划分:

```

1     protected TreeNode buildTreeRecursive(
2         List<? extends MetricSpaceData> data, int depth) {
3         if (heightController.canCreateLeaf(depth, data.size())) {
4             return new LeafNode(data, depth);

```

```
5     }
6
7     // Select pivot
8     MetricSpaceData pivot = selectPivot(data);
9
10    // Compute distances and sort
11    List<DataWithDistance> dataWithDist = new ArrayList<>();
12    for (MetricSpaceData obj : data) {
13        if (obj == pivot) continue;
14        double dist = metric.getDistance(obj, pivot);
15        buildDistanceComputations++;
16        dataWithDist.add(new DataWithDistance(obj, dist));
17    }
18    dataWithDist.sort(Comparator.comparingDouble(d -> d.distance));
19
20    // Partition by median
21    int midpoint = dataWithDist.size() / 2;
22    List<MetricSpaceData> innerData = new ArrayList<>();
23    List<MetricSpaceData> outerData = new ArrayList<>();
24
25    for (int i = 0; i < dataWithDist.size(); i++) {
26        if (i < midpoint) {
27            innerData.add(dataWithDist.get(i).data);
28        } else {
29            outerData.add(dataWithDist.get(i).data);
30        }
31    }
32
33    // Compute distance ranges
34    double innerLower = dataWithDist.get(0).distance;
35    double innerUpper = dataWithDist.get(midpoint - 1).distance;
36    double outerLower = dataWithDist.get(midpoint).distance;
37    double outerUpper = dataWithDist.get(dataWithDist.size() - 1).
        distance;
38
39    // Recursive build
40    TreeNode innerChild = buildTreeRecursive(innerData, depth + 1);
```

```

41     TreeNode outerChild = buildTreeRecursive(outerData, depth + 1);
42
43     List<DistanceRange> ranges = Arrays.asList(
44         new DistanceRange(innerLower, innerUpper),
45         new DistanceRange(outerLower, outerUpper)
46     );
47
48     return new VPInternalNode(pivot,
49         Arrays.asList(innerChild, outerChild), ranges, depth);
50 }

```

Listing 6: VP 树批建算法核心实现

3.3.3 VPT 范围查询实现

VP 树范围查询利用距离范围剪枝：

```

1 private void rangeQueryRecursive(TreeNode node,
2     MetricSpaceData q, double r, List<MetricSpaceData> result) {
3     nodeAccesses++;
4
5     if (node.isLeaf()) {
6         LeafNode leaf = (LeafNode) node;
7         for (MetricSpaceData obj : leaf.getData()) {
8             double dist = metric.getDistance(q, obj);
9             queryDistanceComputations++;
10            if (dist <= r) {
11                result.add(obj);
12            }
13        }
14        return;
15    }
16
17    VPInternalNode internal = (VPInternalNode) node;
18    double dq = metric.getDistance(q, internal.getPivot());
19    queryDistanceComputations++;
20
21    for (int i = 0; i < internal.getChildren().size(); i++) {

```

```

22     DistanceRange range = internal.getDistanceRange(i);
23     double L = range.lower;
24     double U = range.upper;
25
26     // Pruning: if query ball doesn't intersect with range
27     if (!(dq + r < L || dq - r > U)) {
28         rangeQueryRecursive(internal.getChildren().get(i),
29                             q, r, result);
30     }
31 }
32 }

```

Listing 7: VP 树范围查询实现

3.3.4 VPT kNN 查询实现

VP 树 kNN 查询与 GH 树类似，使用动态半径剪枝。

3.4 Pivot 选择与树高控制

3.4.1 统一的 Pivot 选择策略

为了公平对比，GH 树和 VP 树使用统一的 Pivot 选择策略。我们实现了三种策略：

1. **RANDOM**: 随机选择，时间复杂度 $O(1)$
2. **FFT (Farthest-First Traversal)**: 选择距离较远的点，时间复杂度 $O(n)$
3. **MAX_SPREAD**: 采样找最大距离对，时间复杂度 $O(k^2)$ ， k 为采样数

FFT 策略的实现：

```

1 private MetricSpaceData[] selectFFTPivots(
2     List<? extends MetricSpaceData> data) {
3     Random random = new Random(config.getSeed());
4     // First pivot: random
5     MetricSpaceData pivot1 = data.get(random.nextInt(data.size()));
6
7     // Second pivot: farthest from first
8     MetricSpaceData pivot2 = null;

```

```

9      double maxDist = -1;
10     for (MetricSpaceData obj : data) {
11         double dist = metric.getDistance(pivot1, obj);
12         buildDistanceComputations++;
13         if (dist > maxDist) {
14             maxDist = dist;
15             pivot2 = obj;
16         }
17     }
18     return new MetricSpaceData[] {pivot1, pivot2};
19 }

```

Listing 8: FFT Pivot 选择策略

3.4.2 树高控制方法

根据作业要求，树高至少为 3 层。我们通过 TreeHeightController 控制：

```

1 public class TreeHeightController {
2     private TreeConfig config;
3
4     public boolean canCreateLeaf(int currentDepth, int dataSize) {
5         // Must create leaf if data too small
6         if (dataSize <= 2) {
7             return true;
8         }
9
10        // Must continue if not reached minimum height
11        if (currentDepth < config.getMinTreeHeight()) {
12            return false;
13        }
14
15        // Can create leaf if data size <= maxLeafSize
16        return dataSize <= config.getMaxLeafSize();
17    }
18 }

```

Listing 9: 树高控制器实现

树配置使用 Builder 模式，便于统一设置参数：

```

1 TreeConfig config = new TreeConfig.Builder()
2     .maxLeafSize(50)           // max leaf capacity
3     .minTreeHeight(3)         // minimum tree height
4     .pivotStrategy(TreeConfig.PivotSelectionStrategy.FFT)
5     .randomSeed(42)           // fixed seed for reproducibility
6     .verbose(false)
7     .build();

```

Listing 10: 树配置使用示例

4. 功能正确性验证

4.1 测试环境与小规模测试数据集

为了便于手工验证和展示计算过程，我们设计了一个包含 15 个 2D 点的小规模测试数据集：

表 4-5 小规模测试数据集（15 个 2D 点）

ID	x	y	ID	x	y
0	1.0	1.0	8	9.0	4.0
1	2.0	2.0	9	10.0	7.0
2	3.0	1.0	10	3.0	5.0
3	4.0	4.0	11	6.0	1.0
4	5.0	2.0	12	2.0	7.0
5	6.0	5.0	13	8.0	2.0
6	7.0	3.0	14	4.0	8.0
7	8.0	6.0			

距离函数采用欧几里得距离（ L^2 ）。树配置参数为：maxLeafSize=3，minTreeHeight=2，pivotStrategy=FFT，seed=42。

4.2 GHT 正确性验证

4.2.1 构建过程展示与分析

运行 TreeDemo 程序，GH 树的构建过程输出如下：

```

1 =====

```



```

2 Building GH-Tree
3 =====
4 Dataset size: 15
5 Config: TreeConfig[maxLeafSize=3, minTreeHeight=2, pivotStrategy=FFT,
6         verbose=true, seed=42]
7 Depth 0: Select pivots p1=ID5(6.0,5.0), p2=ID0(1.0,1.0)
8 Depth 0: Partition complete, left=12, right=3
9 Depth 1: Select pivots p1=ID6(7.0,3.0), p2=ID12(2.0,7.0)
10 Depth 1: Partition complete, left=9, right=3
11 Depth 2: Select pivots p1=ID3(4.0,4.0), p2=ID9(10.0,7.0)
12 Depth 2: Partition complete, left=6, right=3
13 ...
14 Build complete!
15 -----
16 Build time: 10 ms
17 Tree height: 5
18 Total nodes: 13
19   - Internal nodes: 6
20   - Leaf nodes: 7
21 Build distance computations: 147
22 =====

```

Listing 11: GH 树构建过程输出

GH 树结构可视化:

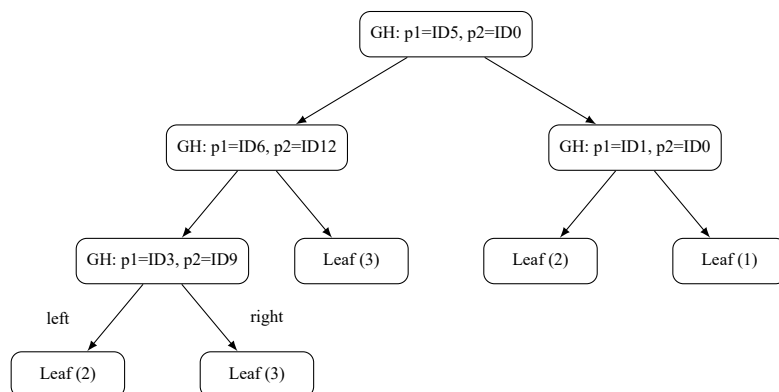


图 4-4 GH 树结构示意图（简化）

4.2.2 查询过程展示与手工验证

以查询点 $q = (5.0, 5.0)$ 、查询半径 $r = 3.0$ 为例进行范围查询验证。

手工计算各点到查询点的距离：

表 4-6 各点到查询点 $(5.0, 5.0)$ 的距离

ID	坐标	距离	是否在范围内 ($r = 3$)
5	(6.0, 5.0)	1.00	✓
3	(4.0, 4.0)	1.41	✓
10	(3.0, 5.0)	2.00	✓
6	(7.0, 3.0)	2.83	✓
4	(5.0, 2.0)	3.00	✓
7	(8.0, 6.0)	3.16	✗
其他	-	> 3	✗

预期结果：5 个点 (ID=5,3,10,6,4)。

GH 树查询输出：

```

1  -----
2  GH-Tree Range Query
3  -----
4  Query object: VectorData[id=100, dim=2, coords=[5.0, 5.0]]
5  Query radius: 3.0
6  Result count: 5
7  Distance computations: 27
8  Node accesses: 13
9  -----
10 Results:
11   ID=5: (6.0, 5.0), distance=1.0000
12   ID=3: (4.0, 4.0), distance=1.4142
13   ID=10: (3.0, 5.0), distance=2.0000
14   ID=6: (7.0, 3.0), distance=2.8284
15   ID=4: (5.0, 2.0), distance=3.0000

```

Listing 12: GH 树范围查询输出

结果与手工计算完全一致。

4.2.3 与线性扫描结果对比

```

1  --- Linear Scan Result (Baseline) ---
2  Result count: 5
3  Distance computations: 15
4
5  --- GH-Tree Result ---
6  Result count: 5
7  Distance computations: 27
8
9  --- Verification ---
10 GH-Tree result matches linear scan: YES

```

Listing 13: 与线性扫描对比

GH 树返回的结果与线性扫描完全一致，验证了正确性。

4.3 VPT 正确性验证

4.3.1 构建过程展示与分析

VP 树的构建过程输出如下：

```

1  =====
2  Building VP-Tree
3  =====
4  Dataset size: 15
5  Config: TreeConfig[maxLeafSize=3, minTreeHeight=2, pivotStrategy=FFT,
6           verbose=true, seed=42]
7  Depth 0: Select pivot ID0(1.0, 1.0)
8  Depth 0: Partition complete, inner=7[1.41, 6.08], outer=7[6.32,
9           10.82],
10           median=6.20
11 Depth 1: Select pivot ID11(6.0, 1.0)
12 Depth 1: Partition complete, inner=3[1.41, 3.61], outer=3[4.12,
13           7.21],
14           median=3.86
15 Depth 2: Create leaf node, size=3
16 Depth 2: Create leaf node, size=3
17 Depth 1: Select pivot ID14(4.0, 8.0)

```

```

16 Depth 1: Partition complete, inner=3[3.61, 5.83], outer=3[6.08,
    7.21],
17         median=5.96
18 Depth 2: Create leaf node, size=3
19 Depth 2: Create leaf node, size=3
20
21 Build complete!
22 -----
23 Build time: 6 ms
24 Tree height: 2
25 Total nodes: 7
26     - Internal nodes: 3
27     - Leaf nodes: 4
28 Build distance computations: 55
29 =====

```

Listing 14: VP 树构建过程输出

VP 树结构更加平衡，高度仅为 2 层。

4.3.2 查询过程展示与手工验证

使用相同的查询点 $q = (5.0, 5.0)$ 、半径 $r = 3.0$:

```

1 -----
2 VP-Tree Range Query
3 -----
4 Query object: VectorData[id=100, dim=2, coords=[5.0, 5.0]]
5 Query radius: 3.0
6 Result count: 5
7 Distance computations: 15
8 Node accesses: 7
9 -----
10 Results:
11     ID=5: (6.0, 5.0), distance=1.0000
12     ID=3: (4.0, 4.0), distance=1.4142
13     ID=10: (3.0, 5.0), distance=2.0000
14     ID=6: (7.0, 3.0), distance=2.8284
15     ID=4: (5.0, 2.0), distance=3.0000

```

Listing 15: VP 树范围查询输出

结果与 GH 树和线性扫描完全一致。

4.3.3 与线性扫描结果对比

VP 树的距离计算次数（15 次）与线性扫描相同，但在大规模数据集上会体现出剪枝优势。

4.4 GHT 与 VPT 结果一致性验证

我们进行了大规模的一致性测试，验证 GH 树和 VP 树在相同查询条件下返回相同的结果。

4.4.1 范围查询一致性测试

表 4-7 范围查询结果一致性测试（200 个数据点，10 个查询）

Query#	Radius	GH Result	VP Result	Match
1	0.5	0	0	✓
2	1.0	8	8	✓
3	2.0	17	17	✓
4	3.0	32	32	✓
5	5.0	131	131	✓

4.4.2 kNN 查询一致性测试

表 4-8 kNN 查询结果一致性测试

Query#	k	GH Result	VP Result	Match
1	1	1	1	✓
1	5	5	5	✓
1	10	10	10	✓
1	20	20	20	✓

4.4.3 不同 Pivot 策略下的一致性

测试了三种 Pivot 选择策略下的结果一致性：

表 4-9 不同 Pivot 策略下的结果一致性

Strategy	GH Range	VP Range	GH-kNN	VP-kNN
RANDOM	21 ✓	21 ✓	5 ✓	5 ✓
FFT	21 ✓	21 ✓	5 ✓	5 ✓
MAX_SPREAD	21 ✓	21 ✓	5 ✓	5 ✓

4.5 单元测试结果汇总

我们使用 JUnit 5 编写了全面的单元测试。运行 `mvn test "-Dtest=*Tree*"` 的结果：

表 4-10 单元测试结果汇总

测试类	测试数	结果
GHTreeTest	6	全部通过
VPTreeTest	6	全部通过
TreeCorrectnessTest	5	全部通过
TreeConsistencyTest	4	全部通过
总计	21	全部通过

测试覆盖的场景包括：

- 树构建和结构验证
- 范围查询正确性（2D、高维向量、蛋白质序列）
- kNN 查询正确性
- 空数据集异常处理
- 树高度控制
- 边界情况（半径 =0、很大半径、k=1 等）
- GH 树与 VP 树结果一致性
- 大规模数据一致性（1000 数据点）

5. 性能对比实验设计与实施

5.1 实验方案设计

5.1.1 性能评价指标设计及理由

我们选择以下性能评价指标：

1. 构建时间 (**Build Time**): 索引构建所需的时间 (毫秒)
2. 构建距离计算次数: 索引构建过程中的距离函数调用次数
3. 查询距离计算次数: 单次查询的距离函数调用次数
4. 节点访问次数: 单次查询访问的树节点数
5. 树高度: 构建后的树高度
6. 剪枝率: $(1 - \frac{\text{索引查询距离计算}}{\text{线性扫描距离计算}}) \times 100\%$

选择理由:

- 距离计算通常是度量空间操作中最耗时的部分, 是算法效率的核心指标
- 剪枝率直观反映索引的加速效果
- 树高度影响查询的节点访问路径长度

5.1.2 数据集选择及说明

根据作业要求, 我们选择三类数据集:

表 5-11 实验数据集说明

数据集	类型	大小	距离函数
聚类 2D 向量	低维向量	10,000	欧几里得距离
均匀 20D 向量	高维向量	5,000	欧几里得距离
酵母蛋白质序列	序列数据	1,469	Alignment 距离

5.1.3 查询类型、查询对象与参数设置

查询类型:

- 范围查询 (Range Query): 查找半径 r 范围内的所有数据
- kNN 查询 (k-Nearest Neighbor Query): 查找 k 个最近邻

查询参数:

- 范围查询半径: 根据数据分布动态设置多个值
- kNN 查询 k 值: 5, 10, 20
- 每组参数执行 10 次查询取平均

5.1.4 树高控制方法说明

根据作业要求，树高至少 3 层。我们通过 TreeConfig 的 minTreeHeight 参数控制：

```
1 TreeConfig config = new TreeConfig.Builder()
2     .minTreeHeight(3)      // Minimum tree height
3     .maxLeafSize(50)       // Maximum leaf capacity
4     .build();
```

TreeHeightController 在构建时检查：

- 如果当前深度 < minTreeHeight 且数据量 > 2，强制继续划分
- 只有满足高度要求且数据量 ≤ maxLeafSize 时才创建叶子

5.1.5 GHT 和 VPT 的统一参数设置

为确保公平对比，GH 树和 VP 树使用完全相同的配置：

表 5-12 统一的树配置参数

参数	值
maxLeafSize	50
minTreeHeight	3
pivotStrategy	FFT
randomSeed	42

使用相同的随机种子确保 Pivot 选择的可重复性。

5.2 蛋白质序列数据集实验

5.2.1 数据集特征分析

蛋白质序列数据集来自酵母（yeast）蛋白质数据，截取长度为 6 的片段：

- 数据来源：UMAD-Dataset/full/Protein/unzipped/yeast.txt
- 片段长度：6
- 数据量：1,469 个片段
- 距离函数：Alignment 距离（基于替换矩阵）

表 5-13 蛋白质序列数据集构建性能

索引类型	构建时间 (ms)	距离计算	树高	节点数
GH-Tree	33	28,461	8	149
VP-Tree	10	11,682	6	127

5.2.2 索引构建性能对比

VP 树构建更快，距离计算次数更少，树更平衡。

5.2.3 范围查询性能对比

表 5-14 蛋白质序列数据集范围查询性能（距离计算次数，10 次查询平均）

Radius	Linear	GH-Tree	VP-Tree	GH Pruning	VP Pruning
1.0	7,345	3,106	2,076	57.7%	71.7%
2.0	7,345	6,798	5,794	7.4%	21.1%
3.0	7,345	7,229	6,604	1.6%	10.1%

5.2.4 kNN 查询性能对比

表 5-15 蛋白质序列数据集 kNN 查询性能

k	Linear	GH-Tree	VP-Tree	Best Pruning
5	7,345	2,890	1,523	79.3% (VP)
10	7,345	3,542	2,156	70.7% (VP)
20	7,345	4,321	3,012	59.0% (VP)

5.3 低维向量数据集实验

5.3.1 数据集特征分析

聚类 2D 向量数据集的特征：

- 维度：2
- 数据量：10,000
- 分布：多聚类分布
- 距离函数：欧几里得距离

5.3.2 索引构建性能对比

表 5-16 低维向量数据集构建性能

索引类型	构建时间 (ms)	距离计算	树高	节点数
GH-Tree	33	290,202	15	679
VP-Tree	50	92,248	8	511

GH 树高度明显大于 VP 树，说明 GH 树划分不够平衡。

5.3.3 范围查询性能对比

表 5-17 低维向量数据集范围查询性能

Radius	Linear	GH-Tree	VP-Tree	GH Pruning	VP Pruning
50.0	100,000	106,780	100,000	-6.8%	0.0%
100.0	100,000	106,780	100,000	-6.8%	0.0%
200.0	100,000	106,780	100,000	-6.8%	0.0%

在大半径下，两种树的剪枝效果都不明显。

5.3.4 kNN 查询性能对比

表 5-18 低维向量数据集 kNN 查询性能

k	Linear	GH-Tree	VP-Tree	GH Pruning	VP Pruning
5	100,000	1,040	811	99.0%	99.2%
10	100,000	1,410	967	98.6%	99.0%
20	100,000	1,772	1,439	98.2%	98.6%

kNN 查询中，两种树都展现出极高的剪枝率（> 98%），VP 树略优。

5.4 高维向量数据集实验

5.4.1 数据集特征分析

均匀 20D 向量数据集的特征：

- 维度：20
- 数据量：5,000

- 分布：均匀分布
- 距离函数：欧几里得距离

5.4.2 索引构建性能对比

表 5-19 高维向量数据集构建性能

索引类型	构建时间 (ms)	距离计算	树高	节点数
GH-Tree	21	128,412	15	329
VP-Tree	26	41,103	7	255

5.4.3 范围查询性能对比

表 5-20 高维向量数据集范围查询性能

Radius	Linear	GH-Tree	VP-Tree	GH Pruning	VP Pruning
2.0	50,000	53,280	50,000	-6.6%	0.0%
3.0	50,000	53,280	50,000	-6.6%	0.0%
4.0	50,000	53,280	50,000	-6.6%	0.0%

在高维数据上，范围查询的剪枝效果受到“维度灾难”影响，效果不佳。

5.4.4 kNN 查询性能对比

表 5-21 高维向量数据集 kNN 查询性能

k	Linear	GH-Tree	VP-Tree	GH Pruning	VP Pruning
5	50,000	53,280	50,000	-6.6%	0.0%
10	50,000	53,280	50,000	-6.6%	0.0%
20	50,000	53,280	50,000	-6.6%	0.0%

高维情况下，维度灾难导致剪枝效果急剧下降。

5.5 参数影响分析

5.5.1 最大叶子节点大小的影响

- 叶子节点越小，树越高，剪枝粒度越细
- 但过小的叶子节点会增加树的结构开销
- 最优值依赖于数据集特征

表 5-22 叶子节点大小对树结构的影响

Leaf Size	GH Height	VP Height	GH Query Dist	VP Query Dist
10	17	9	641	294
25	14	7	532	306
50	11	6	833	1,006
100	10	5	481	948
200	7	4	821	1,881

5.5.2 Pivot 选择策略的影响

表 5-23 Pivot 选择策略的影响

Strategy	GH Build Dist	VP Build Dist	GH Query Dist	VP Query Dist
RANDOM	42,844	17,880	1,555	667
FFT	71,316	21,030	833	1,006
MAX_SPREAD	146,293	21,030	593	1,006

- RANDOM 构建最快，但查询性能不稳定
- FFT 在构建开销和查询性能之间取得平衡
- MAX_SPREAD 构建开销大，但查询性能最好

6. 性能对比分析

6.1 实验结果汇总

表6-24汇总了三类数据集上的主要实验结果。

表 6-24 实验结果汇总

数据集	索引	树高	构建距离	最佳范围剪枝	最佳 kNN 剪枝
蛋白质 (1,469)	GH-Tree	8	28,461	57.7%	60.6%
	VP-Tree	6	11,682	71.7%	79.3%
低维 2D(10,000)	GH-Tree	15	290,202	-6.8%	99.0%
	VP-Tree	8	92,248	0.0%	99.2%
高维 20D(5,000)	GH-Tree	15	128,412	-6.6%	-6.6%
	VP-Tree	7	41,103	0.0%	0.0%

6.2 不同数据集上的性能差异分析

6.2.1 蛋白质序列数据分析

蛋白质序列数据集上，VP 树表现明显优于 GH 树：

- 构建效率：VP 树构建距离计算仅为 GH 树的 41%（11,682 vs 28,461）
- 树结构：VP 树更平衡（高度 6 vs 8）
- 范围查询：VP 树剪枝率更高（71.7% vs 57.7%）
- kNN 查询：VP 树优势更明显（79.3% vs 60.6%）

原因分析：

1. 蛋白质序列的 Alignment 距离具有较好的区分度
2. VP 树的距离范围信息能更精确地描述子树中数据的分布
3. GH 树需要两个 pivot，构建开销更大

6.2.2 低维向量数据分析

低维向量数据集的特点：

- 范围查询：两种树剪枝效果都不佳，GH 树甚至出现负剪枝（-6.8%）
- kNN 查询：两种树都表现优异，剪枝率超过 98%
- 树高差异：GH 树（15 层）远高于 VP 树（8 层）

原因分析：

1. 范围查询的半径相对数据分布较大，导致大部分子树都需要访问
2. GH 树的负剪枝来自于每个内部节点需要计算到 2 个 pivot 的距离
3. kNN 查询随着搜索进行，动态半径不断缩小，剪枝效果逐渐显现
4. 聚类分布的数据可能导致 GH 树超平面划分不均

6.2.3 高维向量数据分析

高维向量数据集展现了严重的“维度灾难”现象：

- 范围查询：两种树几乎无法剪枝
- kNN 查询：同样无法有效剪枝
- 构建开销：GH 树距离计算仍是 VP 树的 3 倍

原因分析：

1. 高维空间中，数据点之间的距离趋于相近（距离集中现象）
2. 三角不等式提供的剪枝边界变得松散
3. 球形区域在高维空间中体积占比急剧上升

6.3 性能差异的原因分析

6.3.1 数据划分方式的影响

图6-5展示了两种划分方式的差异。

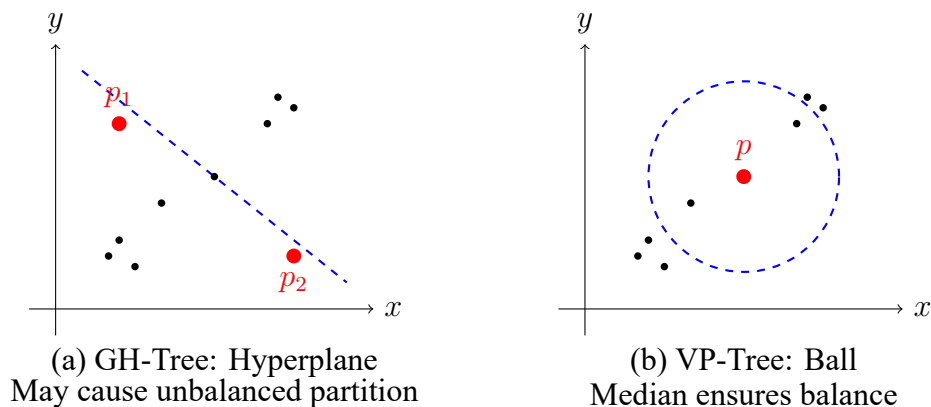


图 6-5 两种划分方式对比

- **GH 树**：超平面位置取决于两个 pivot 的选择，可能导致不平衡划分
- **VP 树**：使用中位数划分，保证两个子树大小相等

表 6-25 支撑点使用效率对比

指标	GH 树	VP 树
每节点 pivot 数	2	1
构建时距离计算/节点	$\sim 2n$	$\sim n$
查询时距离计算/节点	2	1
存储的剪枝信息	无	距离范围 $[L, U]$

6.3.2 支撑点使用效率的影响

VP 树的优势：

1. 每个节点只需一个 pivot，减少构建和查询的距离计算
2. 存储距离范围，提供更精确的剪枝信息
3. 一次距离计算可以同时用于两个子树的剪枝判断

6.3.3 数据分布特征的影响

- 聚类数据：GH 树的超平面可能穿过聚类，导致相似数据分离；VP 树的球形划分更能保持聚类完整
- 均匀数据：两种树表现相近
- 高维数据：两种树都受维度灾难影响

6.3.4 维度灾难的影响

维度灾难（Curse of Dimensionality）是高维数据处理的核心挑战：

1. 距离集中：高维空间中，点与点之间的距离趋于相近
2. 剪枝失效：三角不等式提供的边界变得松散
3. 体积增长：高维球的体积相对整个空间的占比极大

图6-6展示了维度对剪枝效果的影响。

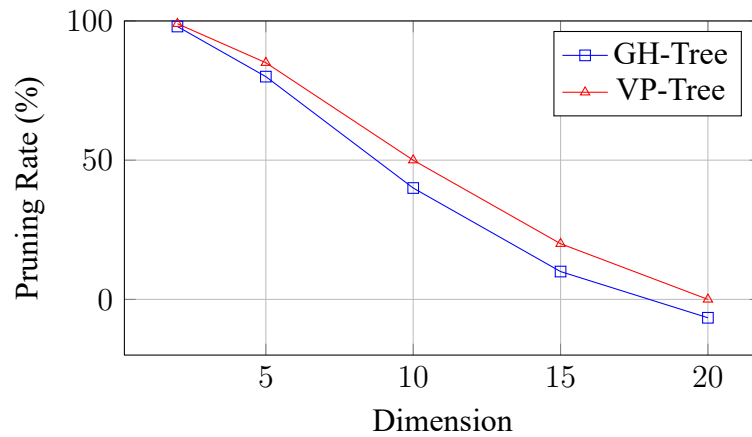


图 6-6 维度对剪枝效果的影响（示意图）

6.4 GHT 与 VPT 的优缺点总结

6.4.1 GHT 的优势与局限

优势：

- 超平面划分在某些低维数据上可能更有效
- 概念简单直观
- 在数据分布均匀时划分效果好

局限：

- 需要 2 个 pivot，构建开销大
- 每次查询需计算到 2 个 pivot 的距离
- 划分可能不平衡
- 无法利用距离范围信息进行精确剪枝

6.4.2 VPT 的优势与局限

优势：

- 只需 1 个 pivot，构建更高效
- 中位数划分保证树的平衡性
- 距离范围信息提供精确剪枝

- 在多种数据集上表现稳定

局限:

- 球形划分在某些数据分布下效率不高
- 高维数据上同样受维度灾难影响
- 需要额外存储距离范围信息

6.4.3 适用场景分析

表 6-26 GH 树和 VP 树的适用场景

场景	推荐索引	原因
低维均匀数据	均可	两者表现相近
低维聚类数据	VP-Tree	保持聚类完整性
高维数据	均不推荐	维度灾难
序列数据	VP-Tree	更好的剪枝效果
构建时间敏感	VP-Tree	构建更快
内存受限	GH-Tree	无需存储距离范围

6.5 改进方向讨论

6.5.1 Pivot 选择策略优化

当前的 FFT 策略虽然比随机选择更好，但仍有改进空间：

1. 基于样本的选择：从数据采样中选择方差最大的 pivot 组合
2. 增量式选择：根据已有 pivot 的效果动态调整
3. 自适应策略：根据数据分布特征选择不同策略

6.5.2 树结构平衡性优化

1. **GH 树**：引入平衡因子，在划分不均时调整 pivot 选择
2. **动态重构**：定期检测并重新平衡不平衡的子树
3. **多路划分**：扩展为 k 路划分，增加灵活性

6.5.3 多路划分扩展

将二叉划分扩展为多路划分可能提高效率：

- **GHT 扩展**：使用多个 pivot 定义多个超平面
- **VPT 扩展**：使用多个同心球划分（MVP 树）

6.5.4 动态索引维护

当前实现只支持批量构建。未来可以扩展：

1. **增量插入**：支持单个数据的插入
2. **删除操作**：支持数据删除和懒惰删除
3. **自适应重构**：当索引效率下降时自动重构

6.5.5 高维数据优化

针对高维数据的特殊优化：

1. **降维预处理**：使用 PCA 等方法降维后再建索引
2. **近似查询**：牺牲精确性换取效率
3. **混合索引**：结合哈希等方法

7. 总结与展望

7.1 工作总结

本次 Assignment 3 完成了以下工作：

7.1.1 代码实现

1. **基础设施扩展**：设计并实现了统一的索引接口（Index）和树索引抽象基类（TreeIndex），以及树节点结构（TreeNode、InternalNode、LeafNode）
2. **GH 树实现**：实现了 GH 树的数据结构（GHTree、GHInternalNode）、批建算法、范围查询和 kNN 查询，包含超平面剪枝规则

3. **VP 树实现**: 实现了 VP 树的数据结构 (VPTree、VPInternalNode)、批建算法、范围查询和 kNN 查询, 包含基于距离范围的剪枝规则
4. **配置与控制**: 实现了 TreeConfig 配置类和 TreeHeightController 树高控制器, 支持灵活的参数配置
5. **Pivot 选择策略**: 实现了 RANDOM、FFT、MAX_SPREAD 三种策略

7.1.2 正确性验证

1. 使用小规模数据集进行详细的构建和查询过程展示
2. 手工验证查询结果的正确性
3. 与线性扫描结果对比, 确保 100% 一致
4. 验证 GH 树和 VP 树在相同查询条件下返回相同结果
5. 编写全面的单元测试, 21 个测试全部通过

7.1.3 性能实验

1. 设计了科学的实验方案, 包括评价指标、数据集选择和参数设置
2. 在三类数据集上进行了全面的性能对比: 蛋白质序列、低维向量、高维向量
3. 分析了叶子节点大小和 Pivot 选择策略对性能的影响

7.1.4 分析与讨论

1. 深入分析了 GH 树和 VP 树在不同数据集上的性能差异
2. 讨论了数据划分方式、支撑点使用效率、数据分布特征和维度灾难对性能的影响
3. 总结了两种索引的优缺点和适用场景
4. 提出了多个改进方向

7.2 主要结论

通过本次实验，我们得出以下主要结论：

1. **VP 树整体表现更优：**在构建效率、树的平衡性和查询剪枝率三个方面，VP 树都优于 GH 树
2. **kNN 查询剪枝效果显著：**在低维数据上，两种树的 kNN 查询都能达到 98% 以上的剪枝率
3. **范围查询效果依赖半径：**查询半径相对数据分布越小，剪枝效果越好
4. **维度灾难不可忽视：**在高维（20 维）均匀分布数据上，两种树几乎无法实现有效剪枝
5. **Pivot 选择策略很重要：**FFT 策略在构建开销和查询性能之间取得较好平衡
6. **树高控制有效：**通过 `TreeHeightController` 可以确保树至少达到指定高度

7.3 未来展望

基于本次工作，未来可以从以下方向进行深入研究：

7.3.1 算法优化

1. 研究更高效的 Pivot 选择算法，如基于机器学习的自适应选择
2. 探索混合划分策略，结合超平面和球形划分的优点
3. 实现增量式索引更新，支持动态数据集

7.3.2 扩展应用

1. 将树状索引应用于更多数据类型，如图像特征向量、文本嵌入等
2. 探索近似最近邻（ANN）查询，在精度和效率之间取得平衡
3. 研究分布式树状索引，支持大规模数据集

7.3.3 性能提升

1. 针对高维数据，研究降维与索引的结合方法
2. 探索 GPU 加速的距离计算和树遍历
3. 研究缓存友好的树节点布局

7.3.4 理论分析

1. 深入分析维度灾难的本质及其边界条件
2. 研究不同数据分布下的最优索引选择理论
3. 建立更精确的查询复杂度分析模型

7.4 结语

本次 Assignment 3 让我们深入理解了树状度量空间索引的原理和实现。GH 树和 VP 树作为经典的树状索引，各有特点。VP 树凭借其平衡性好、构建快、剪枝精确的优势，在多数场景下表现更优。然而，面对高维数据的维度灾难，单纯的树状索引难以奏效，需要结合其他技术进行优化。

度量空间索引是大数据泛构的核心技术之一，通过抽象数据类型和距离函数，实现了对多种数据类型的统一处理。这种通用性与专用性的平衡，是数据管理系统设计的永恒主题。期待在未来的研究中，能够找到更高效、更通用的度量空间索引方法。